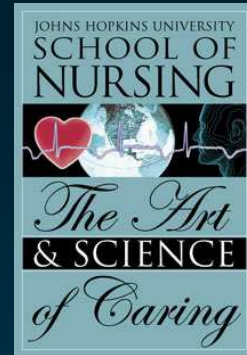Johns Hopkins University
School of Nursing
December 16, 2004

*Introduction to the*
Unified Medical Language System

Olivier Bodenreider

Lister Hill National Center
for Biomedical Communications
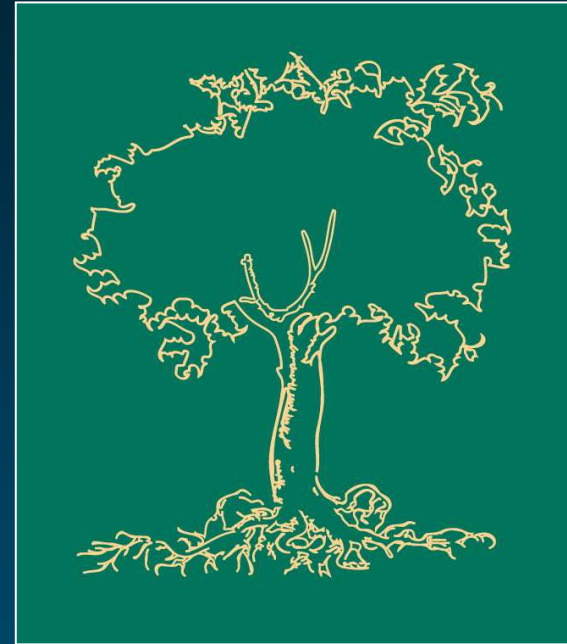Bethesda, Maryland - USA

# Outline

◆ Introduction

◆ Overview through an example

◆ The three UMLS Knowledge Sources

- UMLS Metathesaurus

- UMLS Semantic Network

- SPECIALIST Lexicon and lexical tools

◆ UMLS in action: *MetaMap*

# *Introduction*

# What does **UMLS** stand for?

◆ **U**nified

◆ **M**edical

◆ **L**anguage

◆ **S**ystem

UMLS®
Unified Medical Language System®
UMLS Metathesaurus®

# Motivation

- ◆ Started in 1986

- ◆ National Library of Medicine

- ◆ "Long-term R&D project"

- ◆ Complementary to IAIMS    (Integrated Academic Information Management Systems)

«[…] the UMLS project is an effort to overcome two significant barriers to effective retrieval of machine-readable information.
- The first is the variety of ways the same concepts are expressed in different machine-readable sources and by different people.
- The second is the distribution of useful information among many disparate databases and systems.»

# The UMLS in practice

- ◆ Database
  - Series of relational files
- ◆ Interfaces
  - Web interface: Knowledge Source Server (UMLSKS)
  - Application programming interfaces (Java and XML-based)
- ◆ Applications
  - lvg (lexical programs)
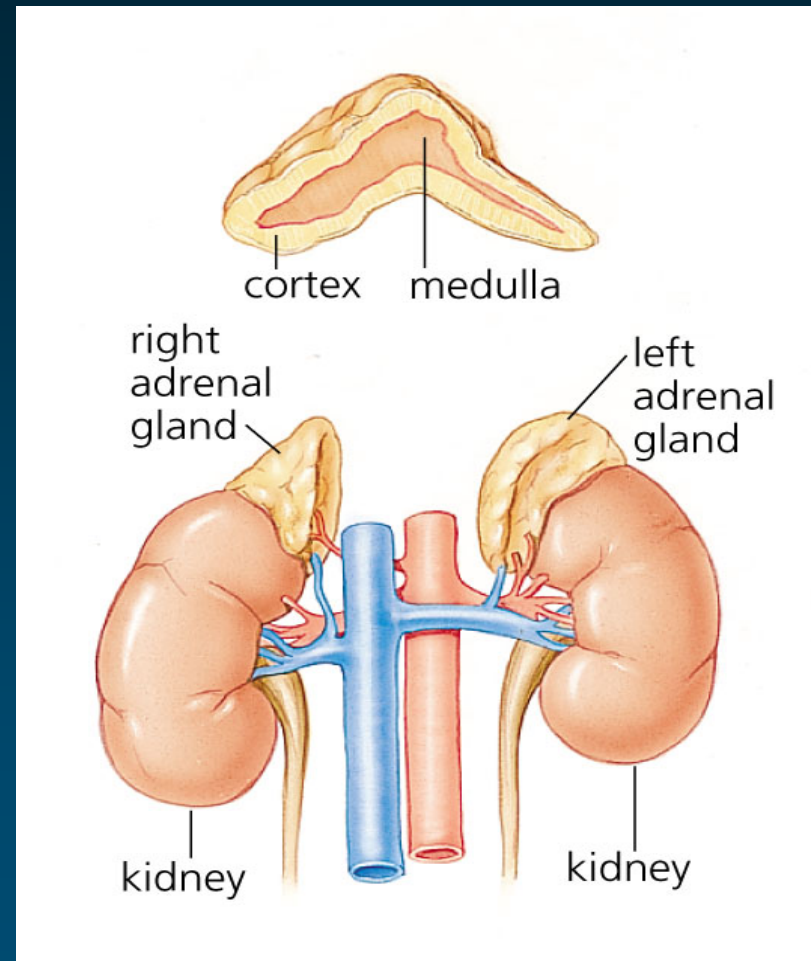  - MetamorphoSys (installation and customization)

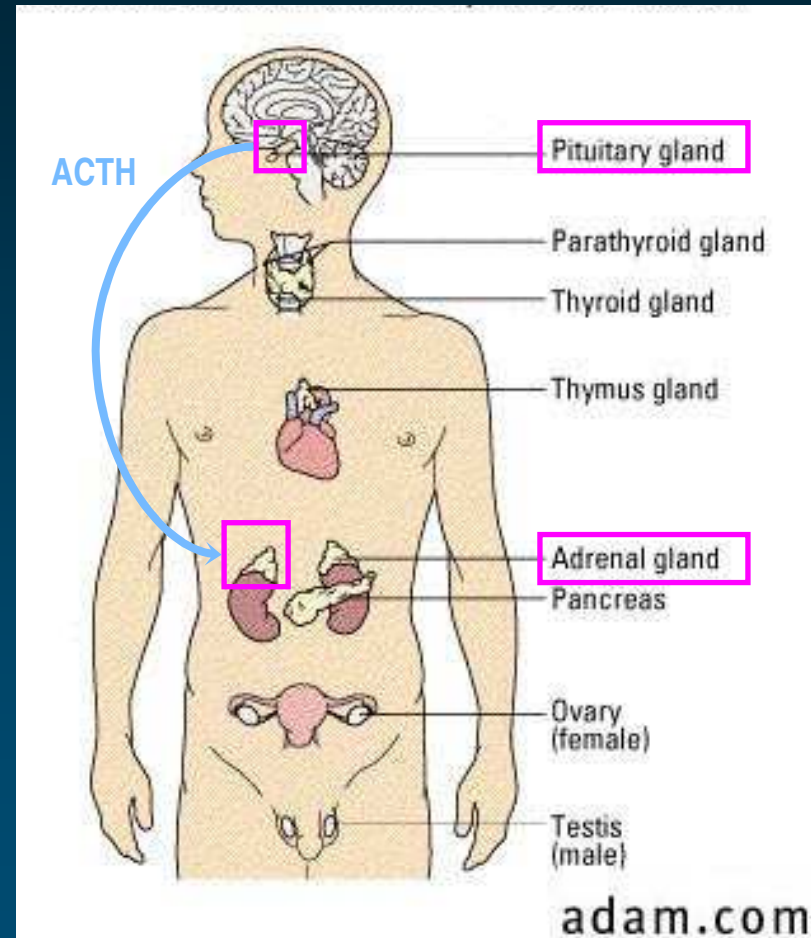The UMLS is *not* an end-user application

6

# *Overview through an example*

# Addison's disease

- Addison's disease is a rare endocrine disorder

- Addison's disease occurs when the adrenal glands do not produce enough of the hormone cortisol

- For this reason, the disease is sometimes called chronic adrenal insufficiency, or hypocortisolism



cortex    medulla

right adrenal gland

left adrenal gland

kidney

kidney

# Adrenal insufficiency  Clinical variants

- ◆ **Primary / Secondary**
  - ● Primary: lesion of the adrenal glands themselves
  - ● Secondary: inadequate secretion of ACTH by the pituitary gland
- ◆ **Acute / Chronic**
- ◆ **Isolated / Polyendocrine deficiency syndrome**



ACTH

Pituitary gland
Parathyroid gland
Thyroid gland
Thymus gland
Adrenal gland
Pancreas
Ovary (female)
Testis (male)

adam.com

# Addison's disease: Symptoms

◆ Fatigue

◆ Weakness

◆ Low blood pressure

◆ Pigmentation of the skin (exposed and non-exposed parts of the body)

◆ …

NLM

# AD in medical vocabularies

- ◆ Synonyms: different terms
  - Addisonian syndrome      ] eponym
  - Bronzed disease
  - Addison melanoderma      symptoms
  - Asthenia pigmentosa
  - Primary adrenal deficiency
  - Primary adrenal insufficiency    clinical
  - Primary adrenocortical insufficiency    variants
  - Chronic adrenocortical insufficiency
- ◆ Contexts: different hierarchies

# Organize terms

- Synonymous terms clustered into a concept
- Preferred term
- Unique identifier (CUI)

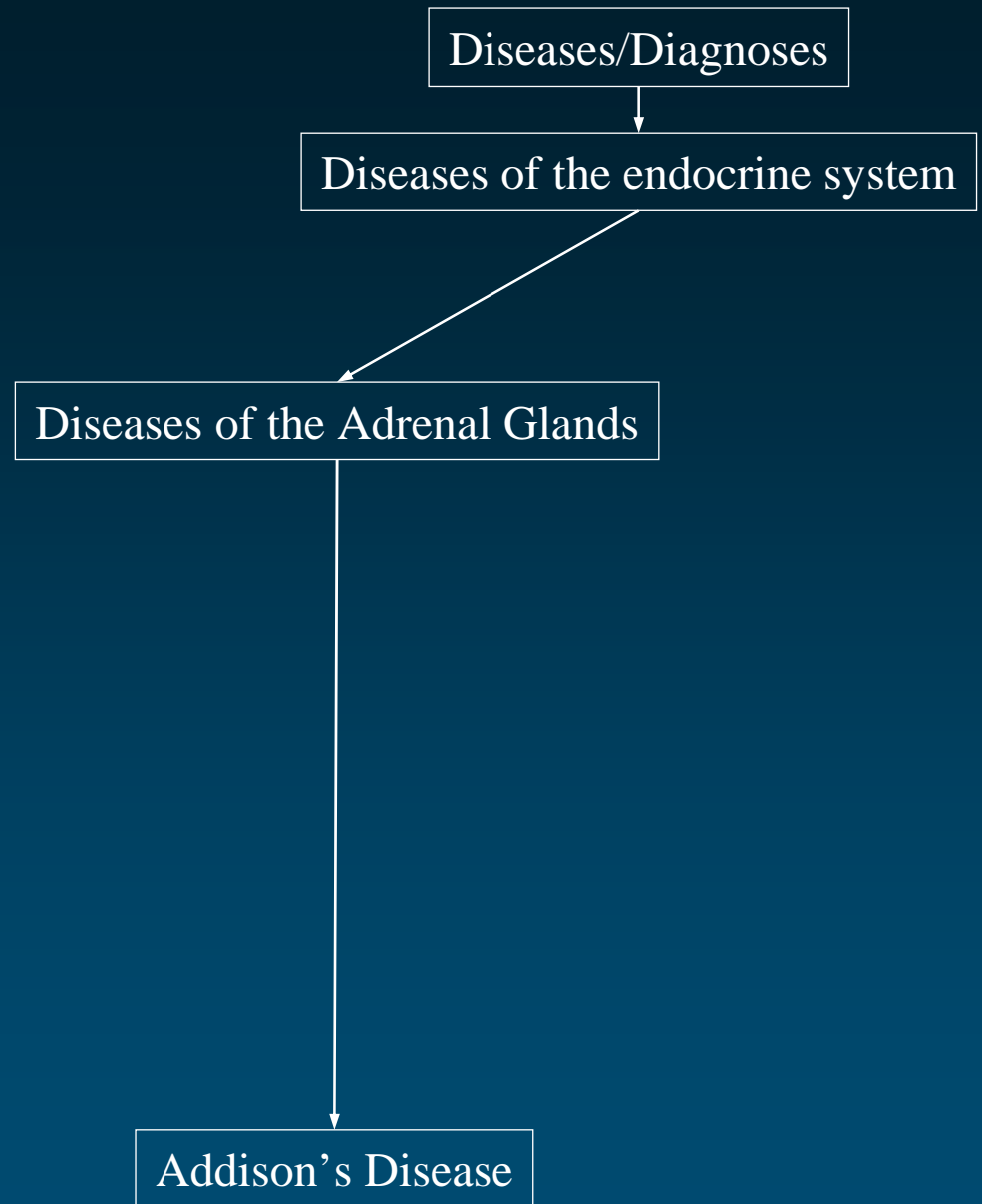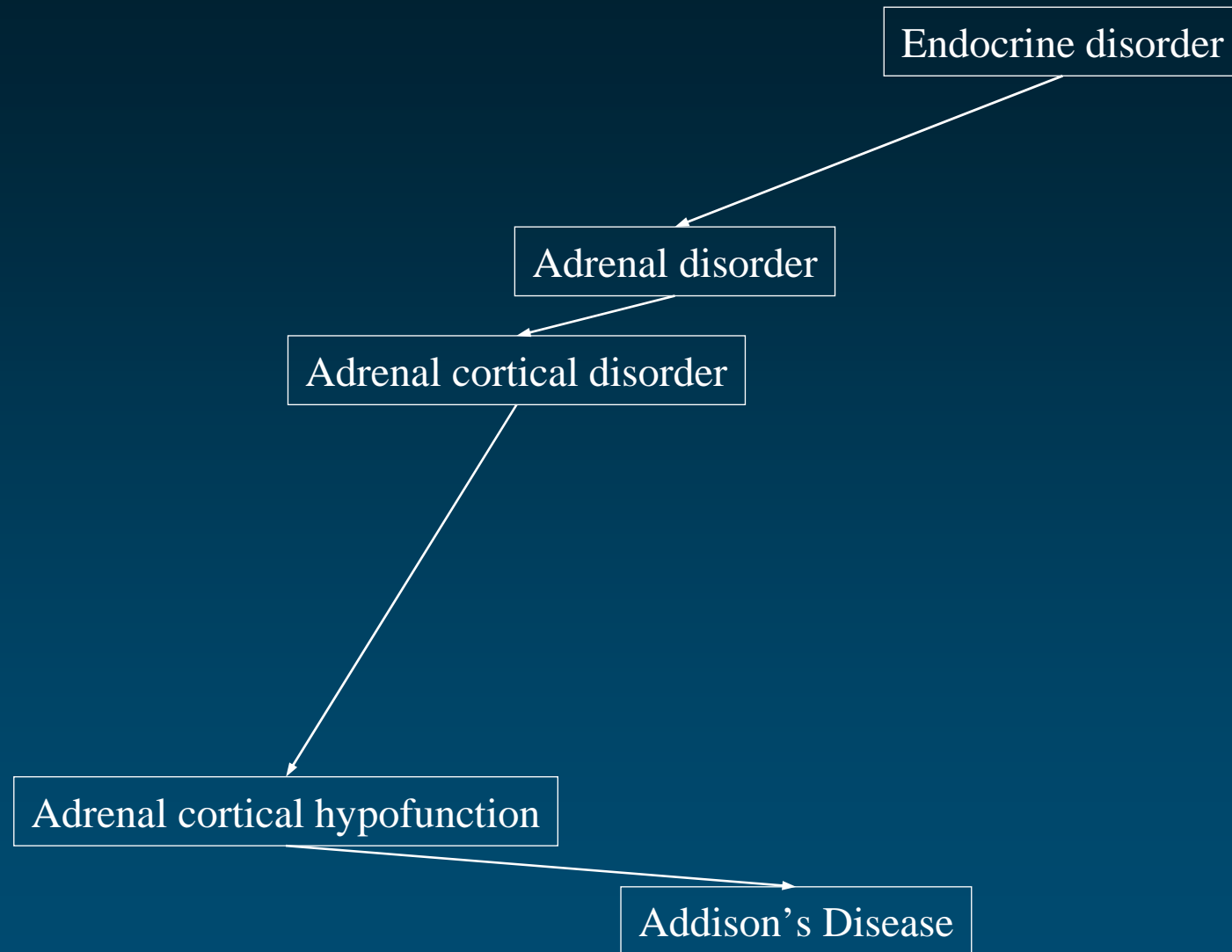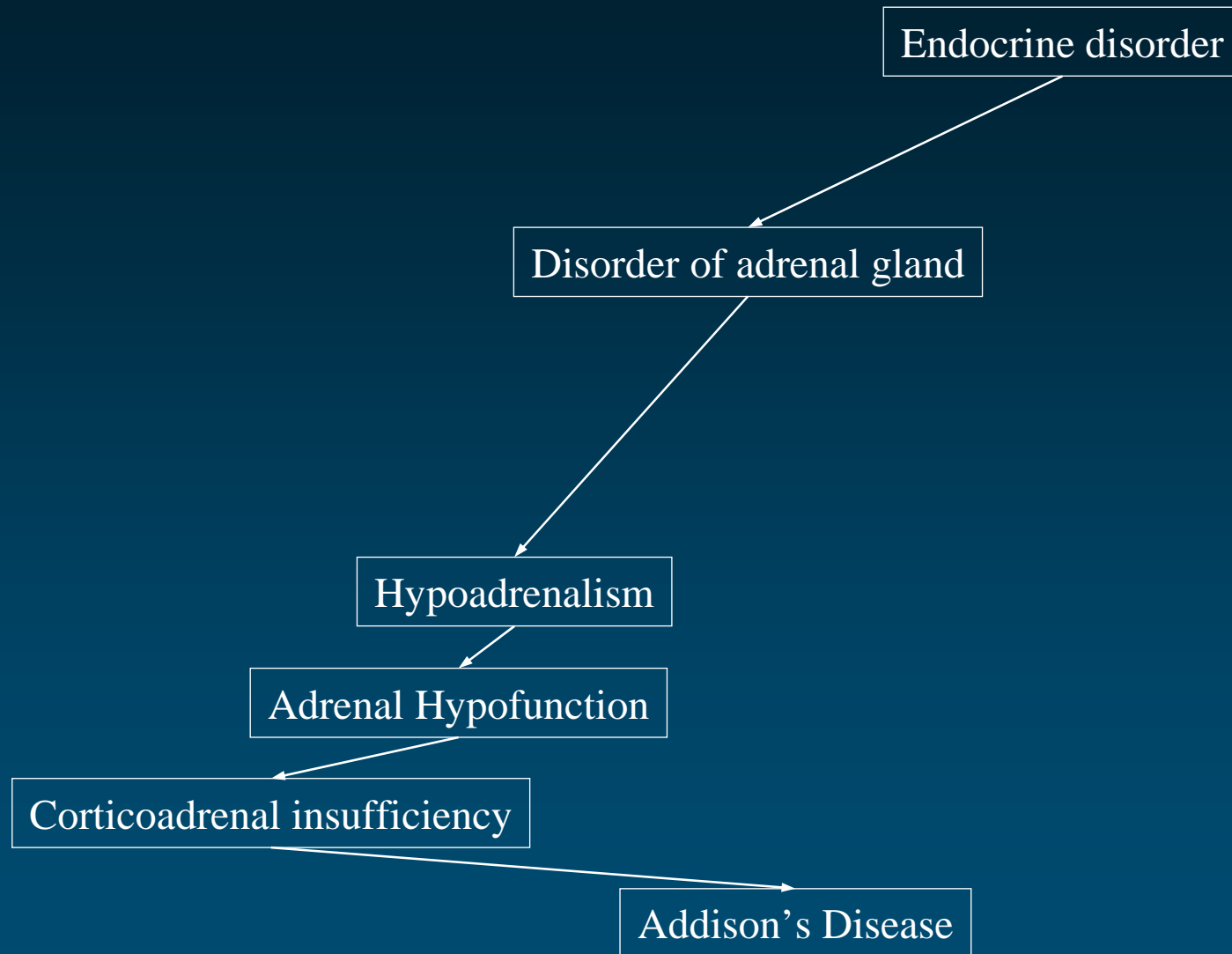| | | |
|---|---|---|
| Adrenal gland diseases | MeSH | D000307 |
| Adrenal disorder | AOD | 0000005418 |
| Disorder of adrenal gland | Read | C15z. |
| Diseases of the adrenal glands | SNOMED | DB-70000 |

C0001621

Adrenal Gland Diseases

**SNOMED International**

Diseases/Diagnoses

Diseases of the endocrine system

Diseases of the Adrenal Glands

Addison's Disease

**MeSH**

Diseases

Endocrine Diseases

Adrenal Gland Diseases

Adrenal Gland Hypofunction

Addison's Disease

**AOD**

Endocrine disorder

Adrenal disorder

Adrenal cortical disorder

Adrenal cortical hypofunction

Addison's Disease
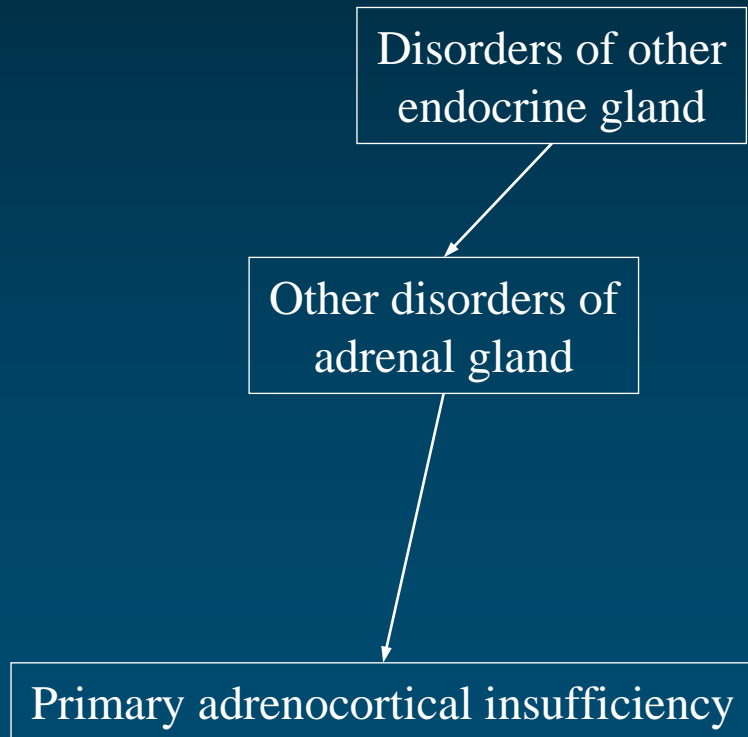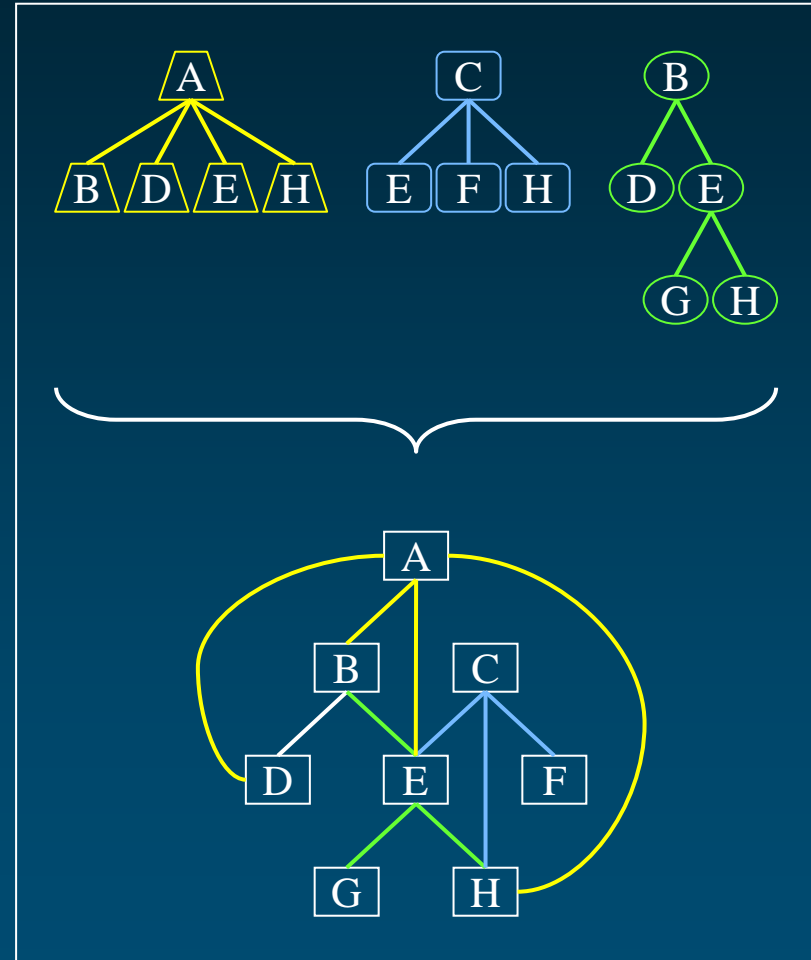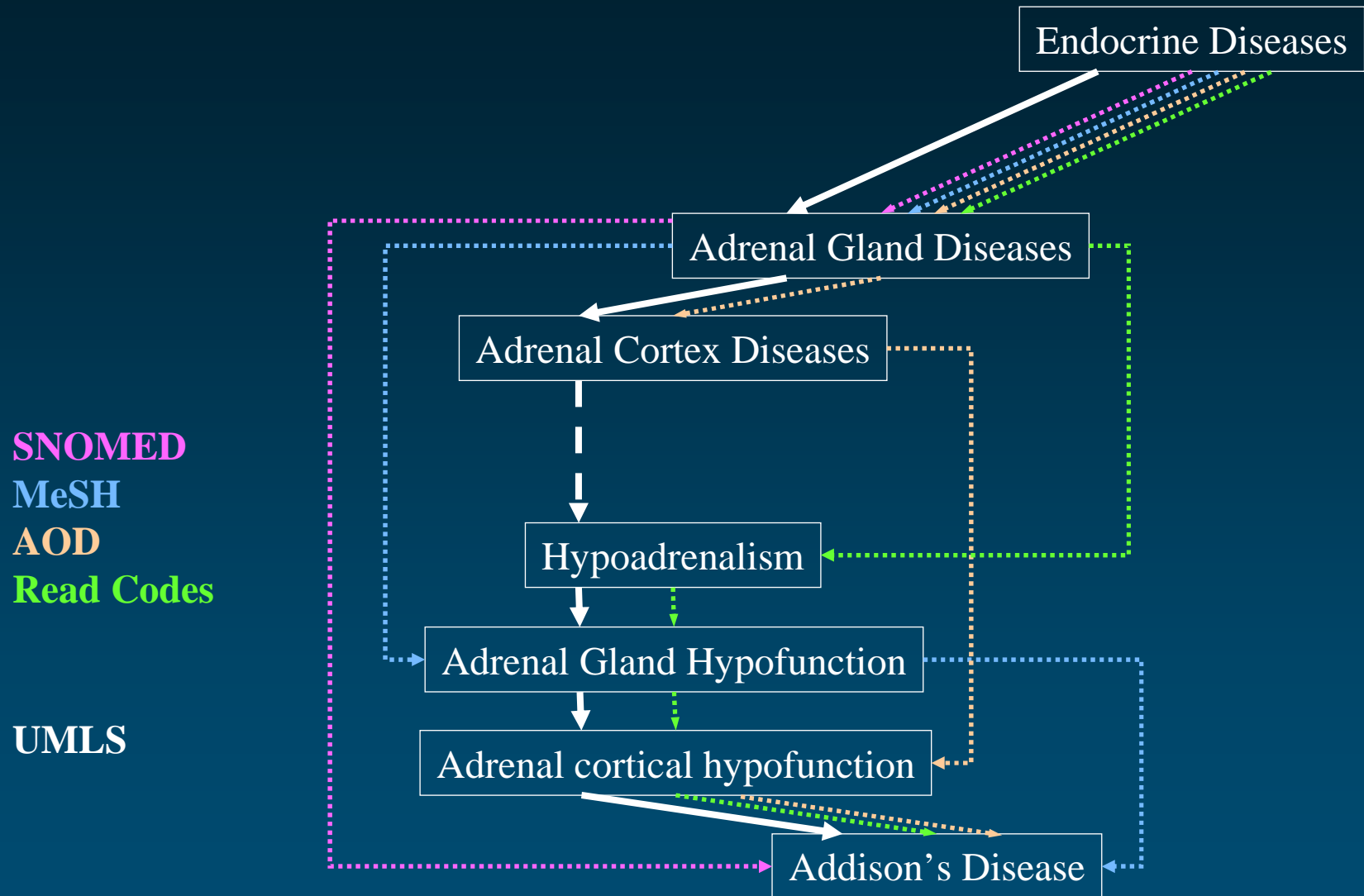
# Organize concepts

- ◆ Inter-concept relationships: hierarchies from the source vocabularies
- ◆ Redundancy: multiple paths
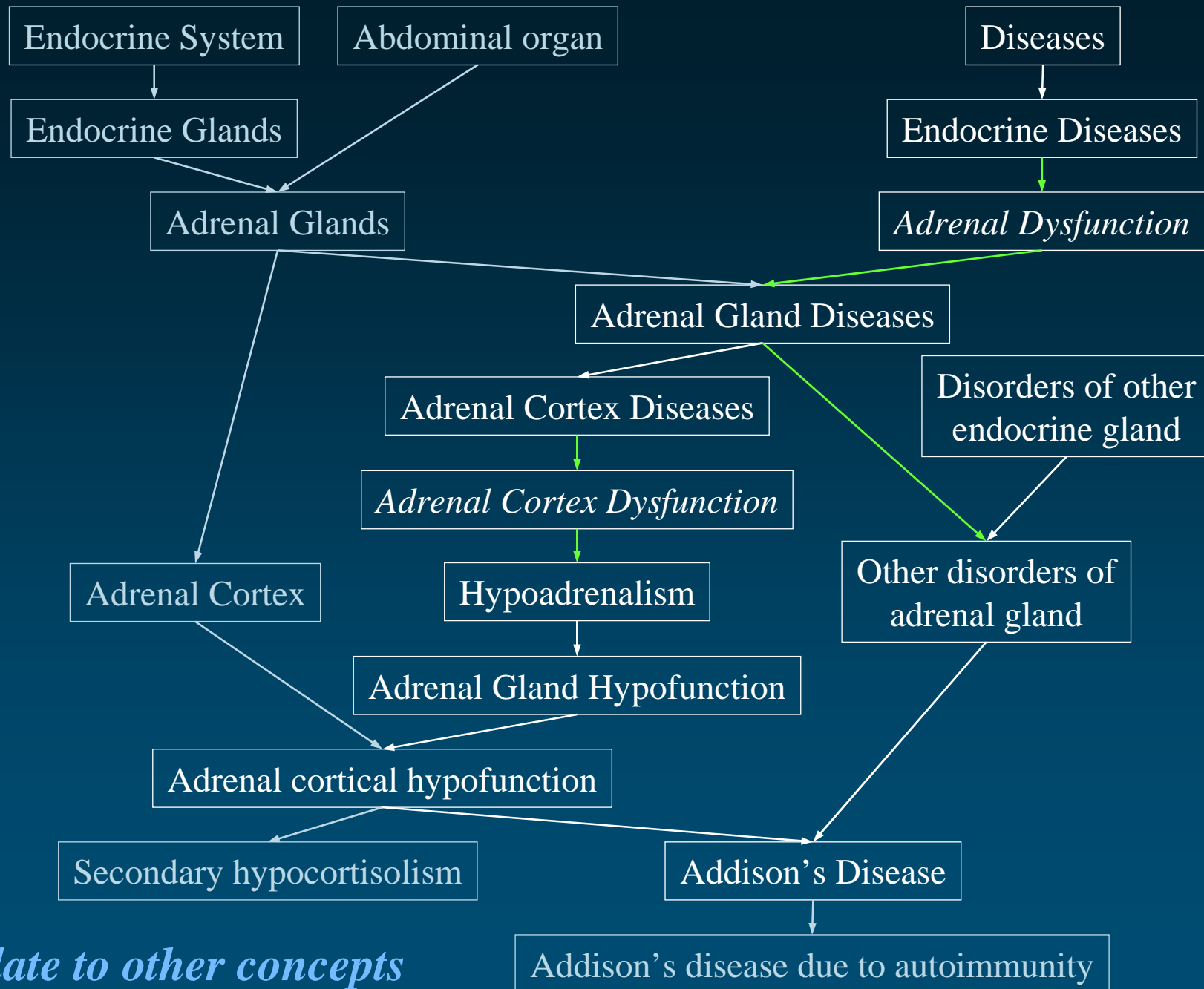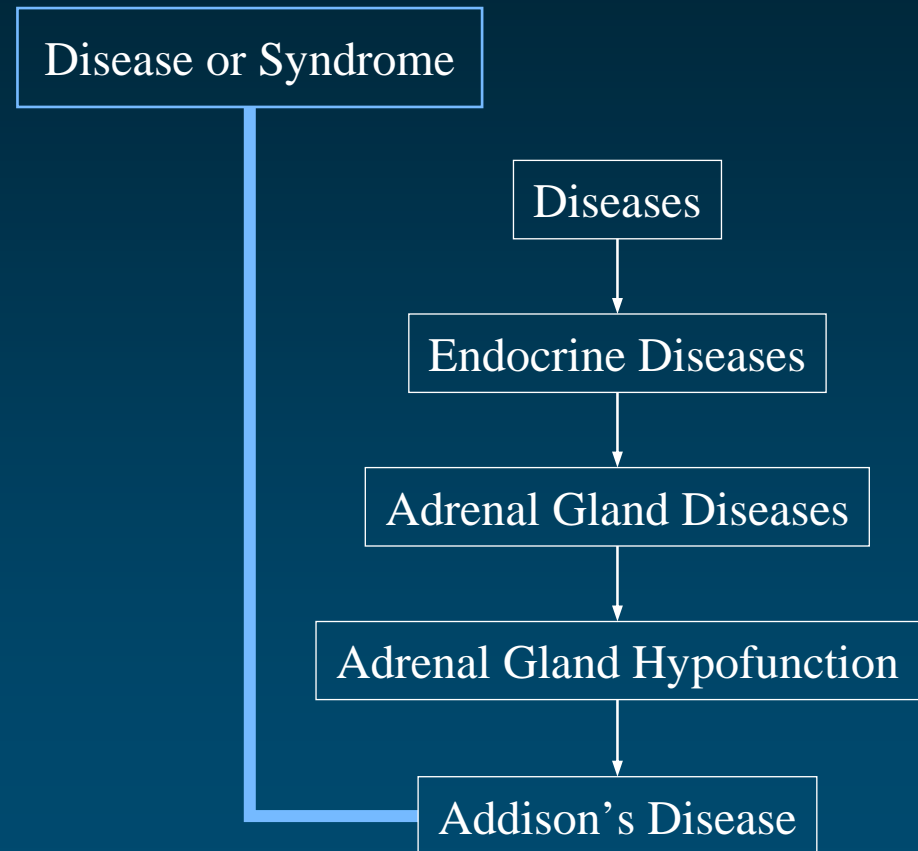- ◆ One graph instead of multiple trees (multiple inheritance)

*organize concepts*

# Relate to other concepts

- Additional hierarchical relationships
  - link to other trees
  - make relationships explicit
- Non-hierarchical relationships
- Co-occurring concepts
- Mapping relationships

# Categorize concepts

- ◆ High-level categories (semantic types)
- ◆ Assigned by the Metathesaurus editors
- ◆ Independently of the hierarchies in which these concepts are located

```
Disease or Syndrome

        Diseases
            │
            ▼
    Endocrine Diseases
            │
            ▼
   Adrenal Gland Diseases
            │
            ▼
 Adrenal Gland Hypofunction
            │
            ▼
     Addison's Disease
```

# How do they do that?

◆ Lexical knowledge

◆ Semantic pre-processing

◆ UMLS editors

# Lexical knowledge

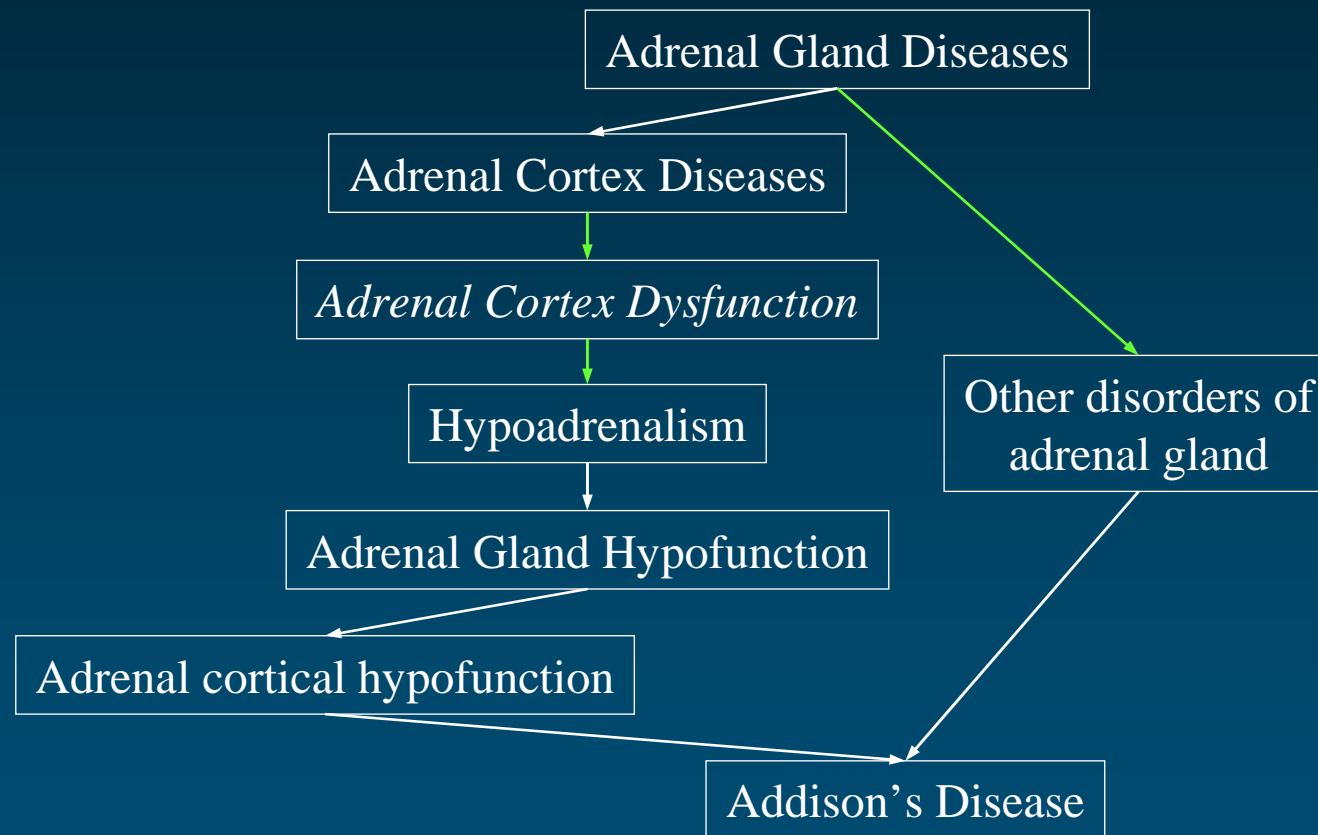Adrenal gland diseases

Adrenal disorder

Disorder of adrenal gland

Diseases of the adrenal glands

C0001621

# Semantic pre-processing

◆ Metadata in the source vocabularies

◆ Tentative categorization

◆ Positive (or negative) evidence for tentative synonymy relations based on lexical features

# Additional knowledge: UMLS editors

Adrenal Gland Diseases

Adrenal Cortex Diseases

*Adrenal Cortex Dysfunction*

Hypoadrenalism

Other disorders of adrenal gland

Adrenal Gland Hypofunction

Adrenal cortical hypofunction

Addison's Disease

# UMLS Summary

◆ Synonymous terms clustered into concepts

◆ Unique identifier


◆ Finer granularity

◆ Broader scope

◆ Additional hierarchical relationships

◆ Semantic categorization

27

# *UMLS Knowledge Sources*

# UMLS  3 components

◆ Metathesaurus

- Concepts
- Inter-concept relationships

◆ Semantic Network

- Semantic types
- Semantic network relationships

◆ Lexical resources

- SPECIALIST Lexicon
- Lexical tools

# UMLS Metathesaurus

# Metathesaurus Basic organization

◆ Concepts

- Synonymous terms are clustered into a concept
- Properties are attached to concepts, e.g.,
  - Unique identifier
  - Definition

◆ Relations

- Concepts are related to other concepts
- Properties are attached to relations, e.g.,
  - Type of relationship
  - Source

# Source Vocabularies

- ◆ 134 source vocabularies
  - ● 126 contributing concept names
- ◆ 73 families of vocabularies
  - ● multiple translations (e.g., MeSH, ICPC, ICD-10)
  - ● variants (American-English equivalents, Australian extension/adaptation)
  - ● subsequent editions usually considered distinct families (ICD: 9-10;  DSM: IIIR-IV)
- ◆ Broad coverage of biomedicine
- ◆ Common presentation

# Biomedical terminologies

- ◆ General vocabularies
  - anatomy (UWDA, Neuronames)
  - drugs (RxNorm, First DataBank, Micromedex)
  - medical devices (UMD, SPN)
- ◆ Several perspectives
  - clinical terms (SNOMED CT)
  - information sciences (MeSH, CRISP)
  - administrative terminologies (ICD-9-CM, CPT-4)
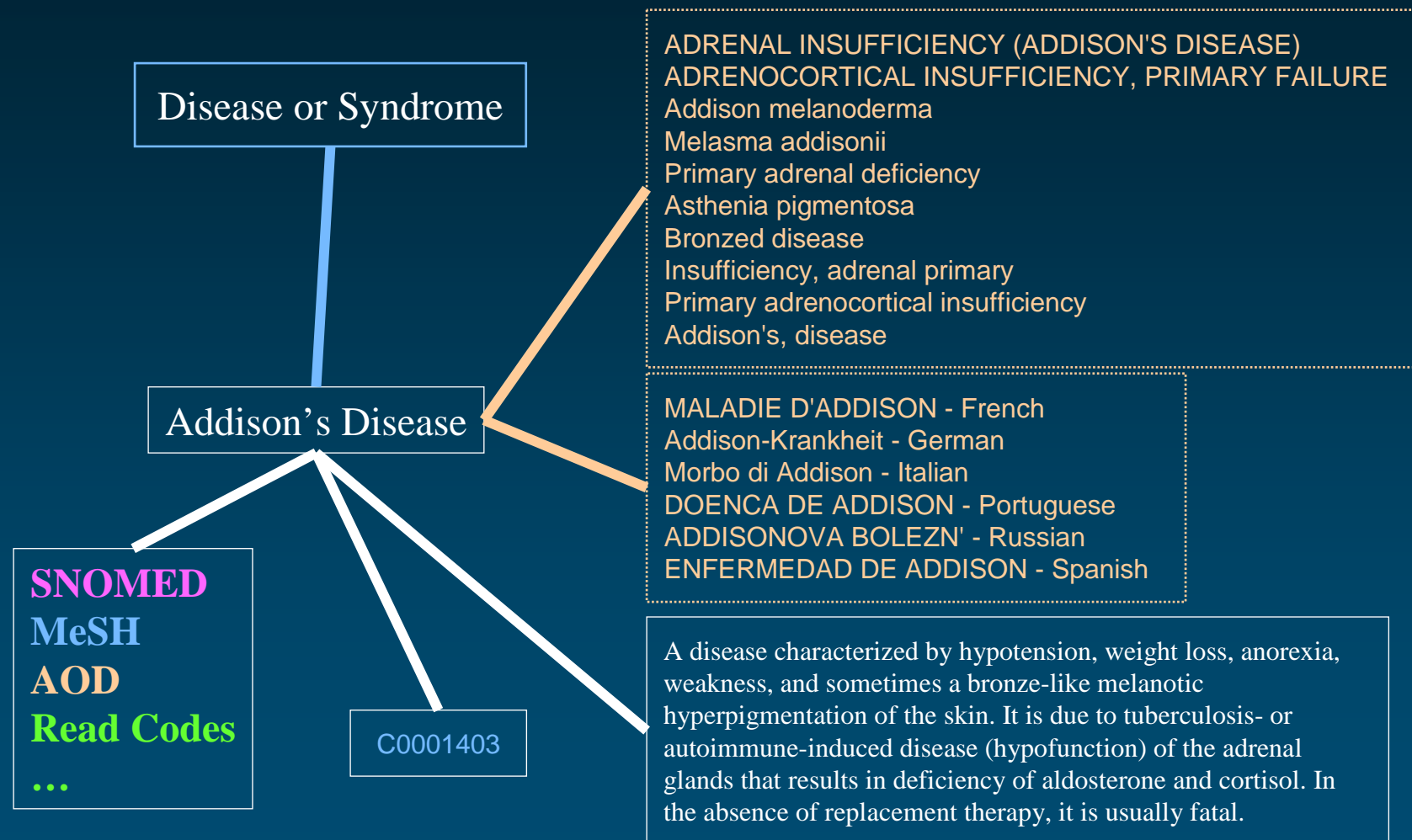  - data exchange terminologies (HL7, LOINC)

# Biomedical terminologies (cont'd)

- **Specialized vocabularies**
  - nursing (NIC, NOC, NANDA, Omaha, PCDS)
  - dentistry (CDT)
  - oncology (PDQ)
  - psychiatry (DSM, APA)
  - adverse reactions (COSTART, WHO ART)
  - primary care (ICPC)
- **Terminology of knowledge bases (AI/Rheum, DXplain, QMR)**

The UMLS serves as a vehicle for the regulatory standards (HIPAA, CHI)

# Addison's Disease: Concept

Disease or Syndrome

Addison's Disease

**SNOMED**
**MeSH**
**AOD**
**Read Codes**
**…**

C0001403

ADRENAL INSUFFICIENCY (ADDISON'S DISEASE)
ADRENOCORTICAL INSUFFICIENCY, PRIMARY FAILURE
Addison melanoderma
Melasma addisonii
Primary adrenal deficiency
Asthenia pigmentosa
Bronzed disease
Insufficiency, adrenal primary
Primary adrenocortical insufficiency
Addison's, disease

MALADIE D'ADDISON - French
Addison-Krankheit - German
Morbo di Addison - Italian
DOENCA DE ADDISON - Portuguese
ADDISONOVA BOLEZN' - Russian
ENFERMEDAD DE ADDISON - Spanish

A disease characterized by hypotension, weight loss, anorexia, weakness, and sometimes a bronze-like melanotic hyperpigmentation of the skin. It is due to tuberculosis- or autoimmune-induced disease (hypofunction) of the adrenal glands that results in deficiency of aldosterone and cortisol. In the absence of replacement therapy, it is usually fatal.

NLM

35

# Metathesaurus Concepts

- ◆ Concept  (> 1M)   CUI
    - ● Set of synonymous concept names
- ◆ Term        (> 3.8 M) LUI
    - ● Set of normalized names
- ◆ String      (> 4.3M)  SUI
    - ● Distinct concept name
- ◆ Atom        (> 5.1M)  AUI
    - ● Concept name in a given source

A0000001  headache   (source 1)
A0000002  headache   (source 2)
S0000001

A0000003  Headache   (source 1)
A0000004  Headache   (source 2)
S0000002

L0000001

A0000005  Cephalgia  (source 1)
S0000003

L0000002

C0000001

# Cluster of synonymous terms

**Concept**
**C0001621**

**Term**
**L0001621**

**S0011232** *Adrenal Gland Diseases*
**S0011231** Adrenal Gland Disease
**S0000441** Disease of adrenal gland
**S0481705** Disease of adrenal gland, NOS
**S0220090** Disease, adrenal gland
**S0044801** Gland Disease, Adrenal

[…]

**Term**
**L0041793**

**S0860744** *Disorder of adrenal gland, unspecified*
**S0217833** Unspecified disorder of adrenal glands

**Term**
**L0161347**

**S0225481** *ADRENAL DISORDER*
**S0627685** DISORDER ADRENAL (NOS)

[…]

**Term**
**L0181041**

**S0632950** *Disorder of adrenal gland*
**S0354509** Adrenal Gland Disorders

[…]

**Term**
**L0368399**

**S0586222** *Adrenal disease*
**S0466921** ADRENAL DISEASE, NOS

[…]

**Term**
**L1279026**

**S1520972** *Nebennierenkrankheiten*

GER

**Term**
**L0162317**

**S0226798** *SURRENALE, MALADIES*

FRE

[…]

# Metathesaurus Evolution over time

- ◆ Concepts never die (in principle)
  - ● CUIs are permanent identifiers
- ◆ What happens when they do die (in reality)?
  - ● Concepts can merge or split
  - ● Resulting in new concepts and deletions

Addison's disease, NOS
C02~~1735~~

Addison's disease
C0001403

1992  1993  1994  1995  1996  1997  1998  1999  …  2004

38

# Metathesaurus Relationships

- Symbolic relations: ~9 M pairs of concepts
- Statistical relations : ~7 M pairs of concepts (co-occurring concepts)
- Mapping relations: 100,000 pairs of concepts

---

- Categorization: Relationships between concepts and semantic types from the Semantic Network

# Symbolic relations

◆ Relation

- Pair of "atom" identifiers
- Type
- Attribute (if any)
- List of sources (for type and attribute)

◆ Semantics of the relationship:
defined by its type [and attribute]

Source transparency: the information
is recorded at the "atom" level

# Symbolic relationships Type

- ◆ Hierarchical
  - ● Parent / Child      **PAR/CHD**
  - ● Broader / Narrower than      **RB/RN**
- ◆ Derived from hierarchies
  - ● Siblings (children of parents)      **SIB**
- ◆ Associative
  - ● Other      **RO**
- ◆ Various flavors of near-synonymy
  - ● Similar      **RL**
  - ● Source asserted synonymy      **SY**
  - ● Possible synonymy      **RQ**

# Symbolic relationships  Attribute

- ◆ Hierarchical
  - ● isa (is-a-kind-of)
  - ● part-of
- ◆ Associative
  - ● location-of
  - ● caused-by
  - ● treats
  - ● …
- ◆ Cross-references (mapping)

# UMLS Semantic Network

# Semantic Network

- **Semantic types (135)**
  - tree structure
  - 2 major hierarchies
    - Entity
      - Physical Object
      - Conceptual Entity
    - Event
      - Activity
      - Phenomenon or Process

# Semantic Network

◆ **Semantic network relationships (54)**

- hierarchical (isa = is a kind of)
    - among types
        - – Animal *isa* Organism
        - – Enzyme *isa* Biologically Active Substance
    - among relations
        - – treats *isa* affects
- non-hierarchical
    - Sign or Symptom *diagnoses* Pathologic Function
    - Pharmacologic Substance *treats* Pathologic Function

# "Biologic Function" hierarchy (isa)

# Associative (non-isa) relationships



48

# Why a semantic network?

◆ Semantic Types serve as high level categories assigned to Metathesaurus concepts, *independently of their position in a hierarchy*

◆ A relationship between 2 Semantic Types (ST) is a possible link between 2 concepts that have been assigned to those STs

- The relationship may or may not hold at the concept level

- Other relationships may apply at the concept level

49

# Relationships can inherit semantics

# SPECIALIST Lexicon
# and lexical tools

# SPECIALIST Lexicon

- Content
  - English lexicon
  - Many words from the biomedical domain
- 200,000+ lexical items
- Word properties
  - morphology
  - orthography
  - syntax
- Used by the lexical tools

# Morphology

- ◆ Inflection
  - ● noun          nucleus, nuclei
  - ● verb          cauterize, cauterizes, cauterized, cauterizing
  - ● adjective      red, redder, reddest
- ◆ Derivation
  - ● verb ⬌ noun      cauterize -- cauterization
  - ● adjective ⬌ noun      red -- redness

# Orthography

◆ Spelling variants

- oe/e                       oesophagus - esophagus

- ae/e                       anaemia - anemia

- ise/ize                   cauterise - cauterize

- genitive mark         Addison's disease
  Addison disease
  Addisons disease

NLM

# Syntax

- ◆ Complementation
  - ● verbs
    - ▪ intransitive
    - ▪ transitive
    - ▪ ditransitive

    I'll treat.

    He treated the patient.

    He treated the patient with a drug.

  - ● nouns
    - ▪ prepositional phrase

      Valve of coronary sinus

- ◆ Position for adjectives

# Lexical tools

- To manage lexical variation in biomedical terminologies
- Major tools
  - Normalization
  - Indexes
  - Lexical Variant Generation program (lvg)
- Based on the SPECIALIST Lexicon
- Used by noun phrase extractors, search engines

# Normalization

| Process | | Text |
|---|---|---|
| | | Hodgkin's diseases, NOS |
| Remove genitive | → | Hodgkin diseases, NOS |
| Remove stop words | → | Hodgkin diseases, |
| Lowercase | → | hodgkin diseases, |
| Strip punctuation | → | hodgkin diseases |
| Uninflect | → | hodgkin disease |
| Sort words | → | disease hodgkin |

# Normalization: Example

Hodgkin Disease
HODGKINS DISEASE
Hodgkin's Disease
Disease, Hodgkin's
Hodgkin's, disease
HODGKIN'S DISEASE
Hodgkin's disease
Hodgkins Disease
Hodgkin's disease NOS
Hodgkin's disease, NOS
Disease, Hodgkins
Diseases, Hodgkins
Hodgkins Diseases
Hodgkins disease
hodgkin's disease
Disease, Hodgkin

normalize → disease hodgkin

# Normalization  Applications

- ◆ Model for lexical resemblance
- ◆ Help find lexical variants for a term
  - Terms that normalize the same usually share the same LUI
- ◆ Help find candidates to synonymy among terms
- ◆ Help map input terms to UMLS concepts

# Indexes

- Word index
  - word to Metathesaurus strings
  - one word index per language
- Normalized word index
  - normalized word to Metathesaurus strings
  - English only
- Normalized string index
  - normalized term to Metathesaurus strings
  - English only

# Lexical Variant Generation program

◆ Tool for specialists (linguists)

◆ Performs atomic lexical transformations

- generating inflectional variants

- lowercase

- …

◆ Performs sequences of atomic transformations

- a specialized sequence of transformations provides the normalized form of a term (the *norm* program)

# MetaMap Motivation

- ◆ Term extraction
  - ● Identifying UMLS concepts from text
- ◆ Usage
  - ● Information indexing and retrieval
  - ● Knowledge extraction / discovery
  - ● Semantic interpretation
- ◆ Characteristics
  - ● Linguistic approach
  - ● Based on UMLS knowledge sources

# MetaMap Methods

◆ Parsing
  - Shallow syntactic analysis
  - SPECIALIST lexicon
  - Xerox part-of-speech tagger

◆ Variant generation

◆ Candidate retrieval
  - Retrieve candidate terms containing at least one variant

◆ Candidate evaluation
  - Rank candidate terms with respect to closeness to input text (centrality, variation, coverage, and cohesiveness)

# MetaMap Example

Molluscum Contagiosum
C0026393

Disease
C0012634

causes
C0015127

Pox virus (Poxviridae)
C0032868

Causing
C0678227

Causation
C0085978

Molluscum contagiosum is a disease caused by a poxvirus of the Molluscipox virus genus that produces a benign self-limited papular eruption of multiple umbilicated cutaneous tumors.

Virus
C0042776

C0260037 Multiple tumors

C0037286 Cutaneous tumor

C0221912 Cutaneous

C0037267 Skin

Papular eruption C0221202

Cutaneous eruption C0332474

Benign C0205183

Papular C0332564

65 […]

# Using MetaMap  MMTx

- Requires UMLS license
- Local implementation (Java-based)
- Provides
  - Stand-alone application
  - API for integrating in other applications

http://mmtx.nlm.nih.gov

# Medical Ontology Research

Contact: olivier@nlm.nih.gov
Web: mor.nlm.nih.gov



*Olivier Bodenreider*

Lister Hill National Center
for Biomedical Communications
Bethesda, Maryland - USA

# Appendix

# Knowledge Source Server
*Web Interface*

http://umlsks.nlm.nih.gov

# UMLS Knowledge Source Server Home Page

{UMLS_2002} UMLS® Semantic Navigator ¤ [2.07] - Netscape 6

**Siblings**

**Concepts & Ideas**

- Clinical Syndromes ¤

**Disorders**

- Acquired Immunodeficiency Syndrome ¤
- Acute adrenal insufficiency ¤
- Addisonian crisis ¤
- Adrenal atrophy ¤
- Adrenal calcification ¤
- Adrenal hemorrhage ¤
- Adrenal infarction ¤
- Adrenal insufficiency due to adrenal metastasis ¤
- Adrenogenital Syndrome ¤
- Allergic arthritis ¤
- Angelman Syndrome ¤
- Asperger syndrome <1> ¤
- Autoerythrocyte sensitivity

Hypoadrenalism

Collagen Diseases

DISEASES OF THE IMMUNE SYSTEM: GENERAL TERMS

Adrenal Gland Diseases

ypofunction

Autoimmune Diseases

Other disorders of adrenal gland

Addison's Disease

ddison's disease with drenoleucodystrophy

Addison's disease due to autoimmunity

Addisons Disease Secondary To Idiopathic Atrophy

Addiso

**Other Related Concepts**

**Disorders**

- Addisonian crisis ¤
- Autoimmune Syndrome Type II, Polyglandular ¤
- Tuberculosis ¤
- Tuberculosis of adrenal glands ¤
- Tuberculous Addison's disease ¤

(5 other related

**Co-occurring Concepts**

**Anatomy**

- Adrenal Cortex [14] ¤
- Adrenal Glands [17] ¤
- Liver [2] ¤
- Tears body substance [2] ¤
- X Chromosome [3]

**Chemicals & Drugs**

BCI     **Addison's Disease**     LEGEND *

Start again    Apply new parameters

Restrict to vocabulary:   Show all

Highlight vocabulary:   Nothing

UMLS data:   UMLS_2002

Type of hierarchical   ⦿ All ○ Parent/Child only ○

**Similar Concepts**

- Adrenal cortical hypofunction ¤

(1 concept)

**Closest MeSH Terms**

**Main Headings**

- Addison's Disease

Document: Done (2.5 secs)

# Knowledge Source Server
## *Application Programming Interface*

# UMLSKS API basics

- ◆ Remote server at NLM
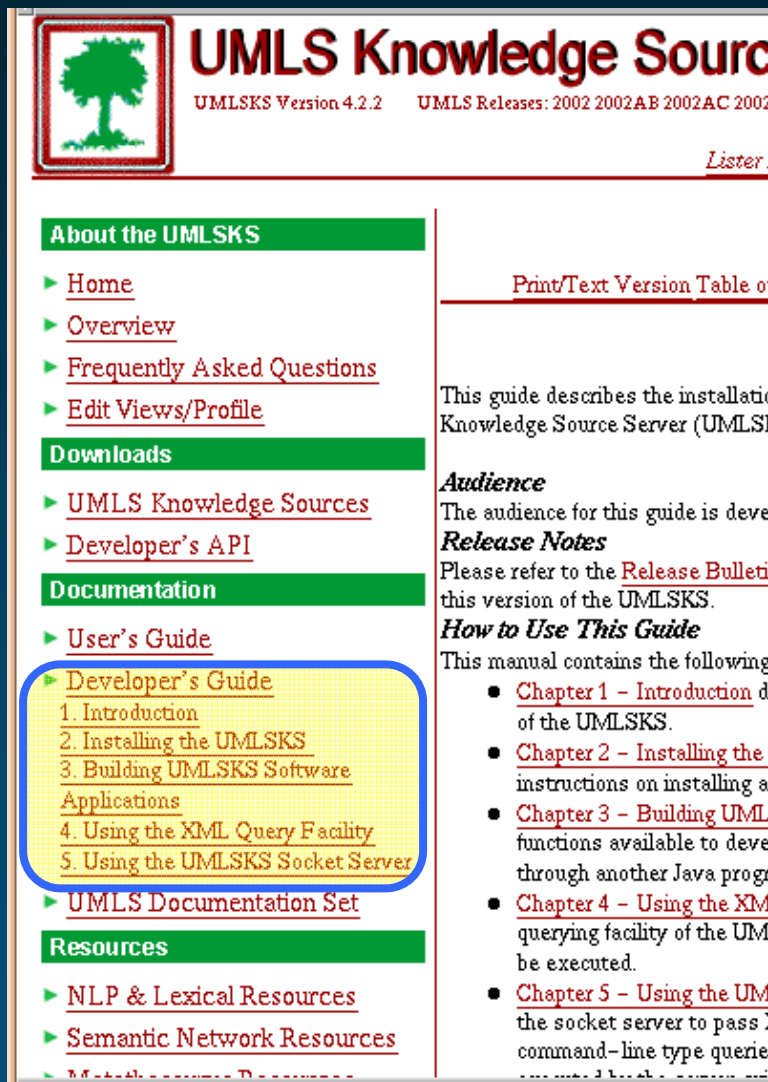- ◆ Local application connected through

| Java RMI | TCP/IP socket |
|---|---|
| ◆ Java-based applications | ◆ XML-based queries |
| ◆ Developer's Guide: Chapter 3 | ◆ Developer's Guide: Chapter 5 |
| ◆ Set of Java classes (part of the UMLSKS API download) | ◆ XML schema |
| ◆ Detailed *Javadoc* documentation online and with API download | ◆ Socket server<br> ● Host: umlsks.nlm.nih.gov<br> ● Port: 8042 |

74

# Developer's Guide

# MetamorphoSys

# What is MetamorphoSys?

◆ Tool distributed with the UMLS

◆ Multi-platform Java software

◆ The UMLS installation and customization wizard
- Installs Knowledge Sources to local storage
- Subsets and customizes a local Metathesaurus

# Why use MetamorphoSys?

*Customize the Metathesaurus*

◆ To remove terminology that is unhelpful, or even harmful, to your needs and purposes

◆ To comply with terms of license agreement

*Changing Default Settings*

◆ To alter the preferred name

◆ To alter suppressibility of specific source term types

# Bibliography

# UMLS documentation and support

◆ UMLS homepage        http://umlsinfo.nlm.nih.gov/

  ● with links to all other UMLS information

◆ UMLSKS homepage      http://umlsks.nlm.nih.gov/

  ● with links to the User's and Developer's guides

◆ Email address for support    custserv@nlm.nih.gov

# References

- ◆ UMLS as a research project
  - ● Lindberg, D. A., Humphreys, B. L., & McCray, A. T. (1993). The Unified Medical Language System. *Methods Inf Med, 32*(4), 281-91.
  - ● Humphreys, B. L., Lindberg, D. A., Schoolman, H. M., & Barnett, G. O. (1998). The Unified Medical Language System: an informatics research collaboration. *J Am Med Inform Assoc, 5*(1), 1-11.
  - ● Bodenreider O. (2004). The Unified Medical Language System (UMLS): Integrating biomedical terminology. *Nucleic Acids Research*; D267-D270.

81

# References

- Technical papers
  - McCray, A. T., & Nelson, S. J. (1995). The representation of meaning in the UMLS. *Methods Inf Med, 34*(1-2), 193-201.
  - Campbell, K. E., Oliver, D. E., Spackman, K. A., & Shortliffe, E. H. (1998). Representing thoughts, words, and things in the UMLS. *J Am Med Inform Assoc, 5*(5), 421-31.
- Comprehensive bibliography 1986-96

  http://www.nlm.nih.gov/pubs/cbm/umlscbm.html